

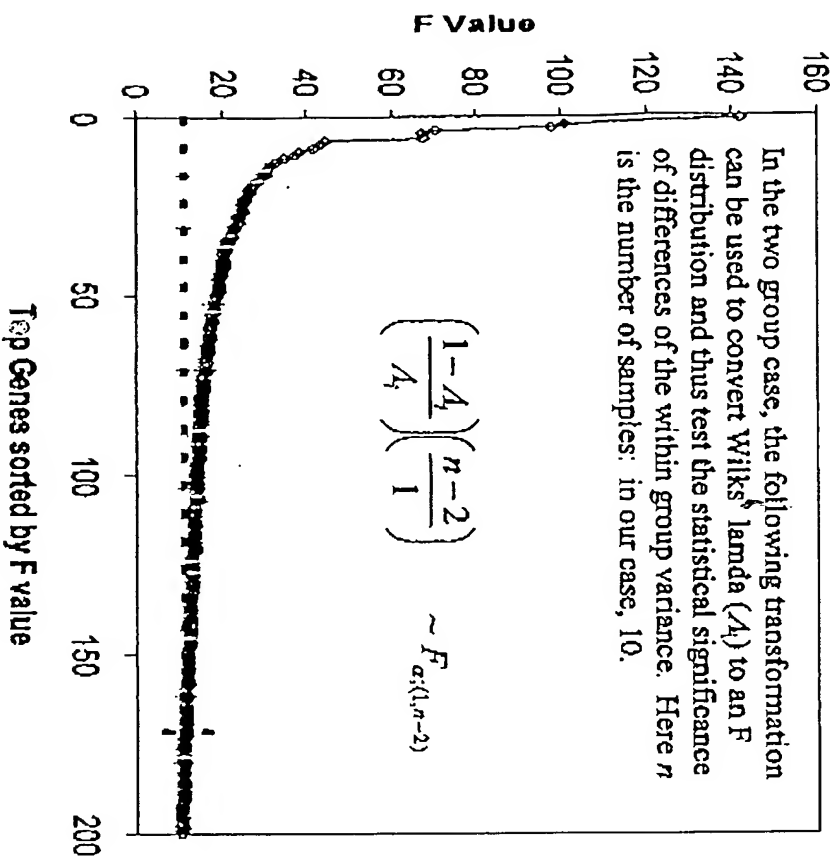
**Figure 1:** A comparison of two idealized gene expression distributions. Gene A represents a good discriminator, because the variance within each group is significantly smaller than the between-group variance. Gene B, on the other hand, has significant overlap between the two groups relative to their variances. To capture this distinction for a gene  $i$  that may be classified in  $j$  classes, we calculate the total variance ( $V_i^t$ ) and compare it to the within cluster variance ( $V_i^c$ ) to determine the Wilks' lambda score ( $\Lambda_i$ ): lower scores mean better discrimination.

$$V_i^t = (\mathbf{x}_i - \bar{\mathbf{x}}_i)^T (\mathbf{x}_i - \bar{\mathbf{x}}_i)$$

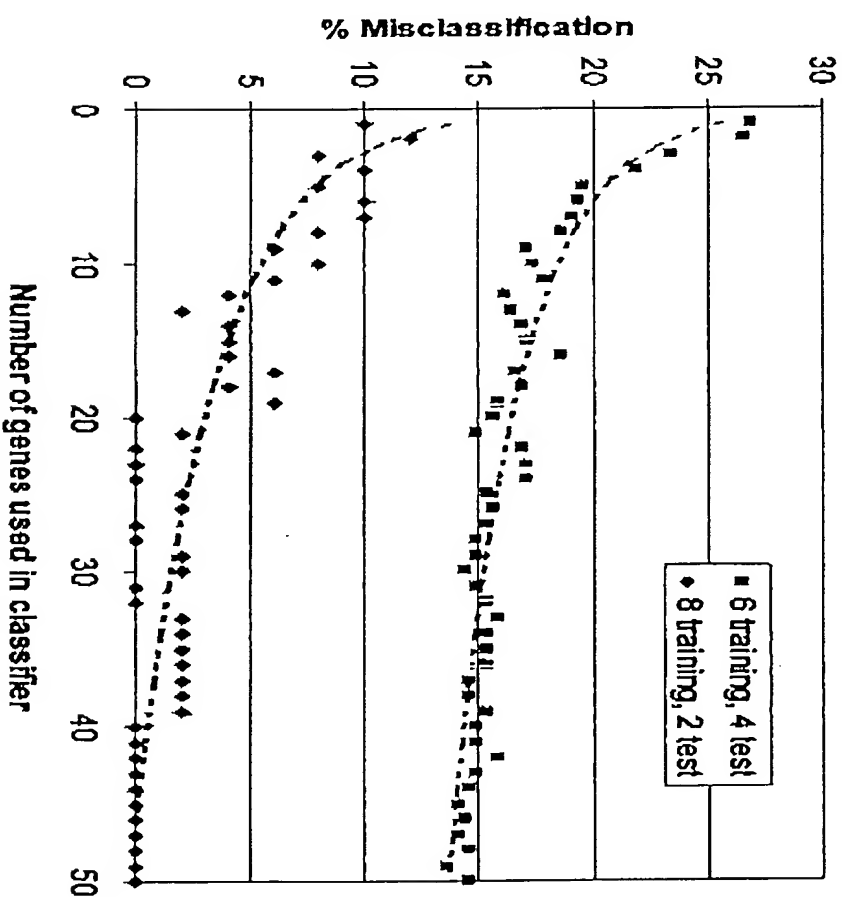
$$V_i^c = \sum_{j=1}^n V_i^{c,j} = \sum_{j=1}^n (\mathbf{x}_i^j - I\bar{\mathbf{x}}_i^j)^T (\mathbf{x}_i^j - I\bar{\mathbf{x}}_i^j)$$

$$\Lambda_i = \frac{V_i^c}{V_i^t}$$

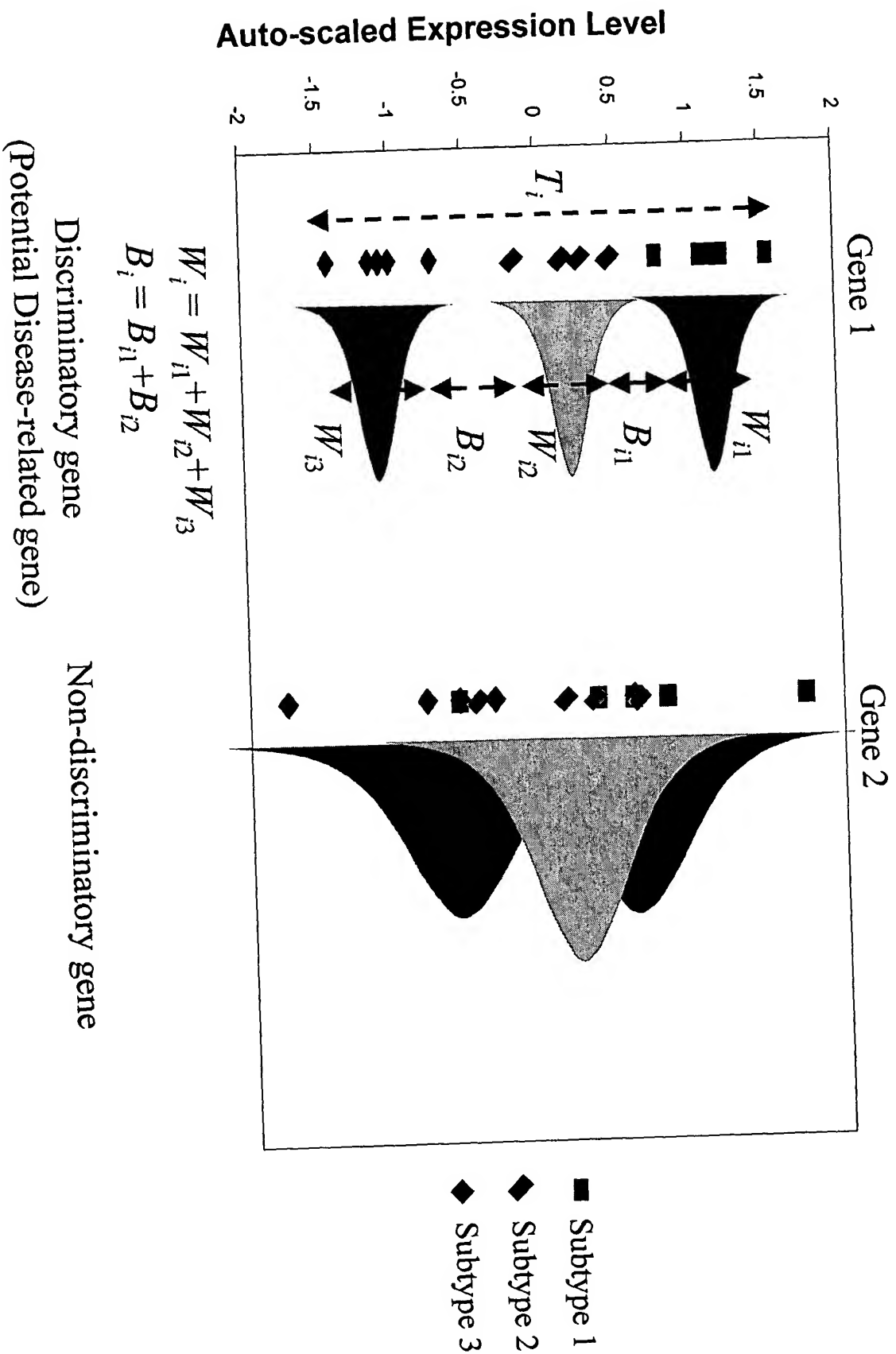
### a) Approximation of Wilks's lambda by a univariate F statistic



### b) Classification Error Rates



**Figure 2:** a) The top 200 discriminating genes rank ordered by the F statistic (equation in the figure). The top 171 genes give an F value above the  $\alpha = 0.99$  confidence interval ( $F_{\alpha; (1,9)} = 11.26$ ), providing a conservative estimate of the genes characteristic of the cancerous state. b) Cross-validation results of classification through the CDA analysis to determine discriminatory genes. In each case, the best discriminatory genes were determined independently of the other cases, so each classification result arises from a different sub-set of the genes. Perfect classification requires about 45 genes when eight samples are included in the training set.



**Figure 3.** Discriminatory genes (potential disease related genes) and Non-discriminatory genes.

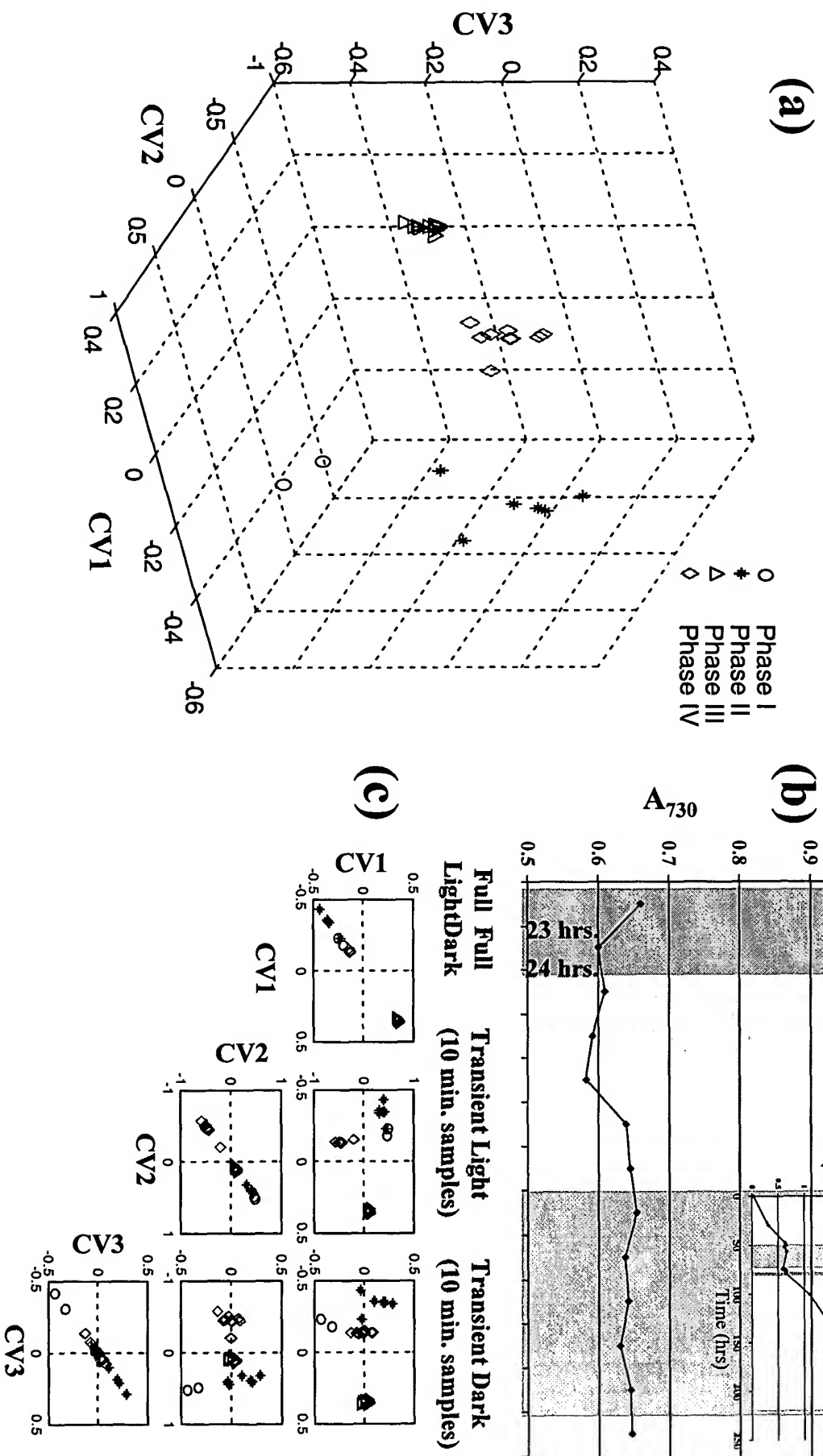
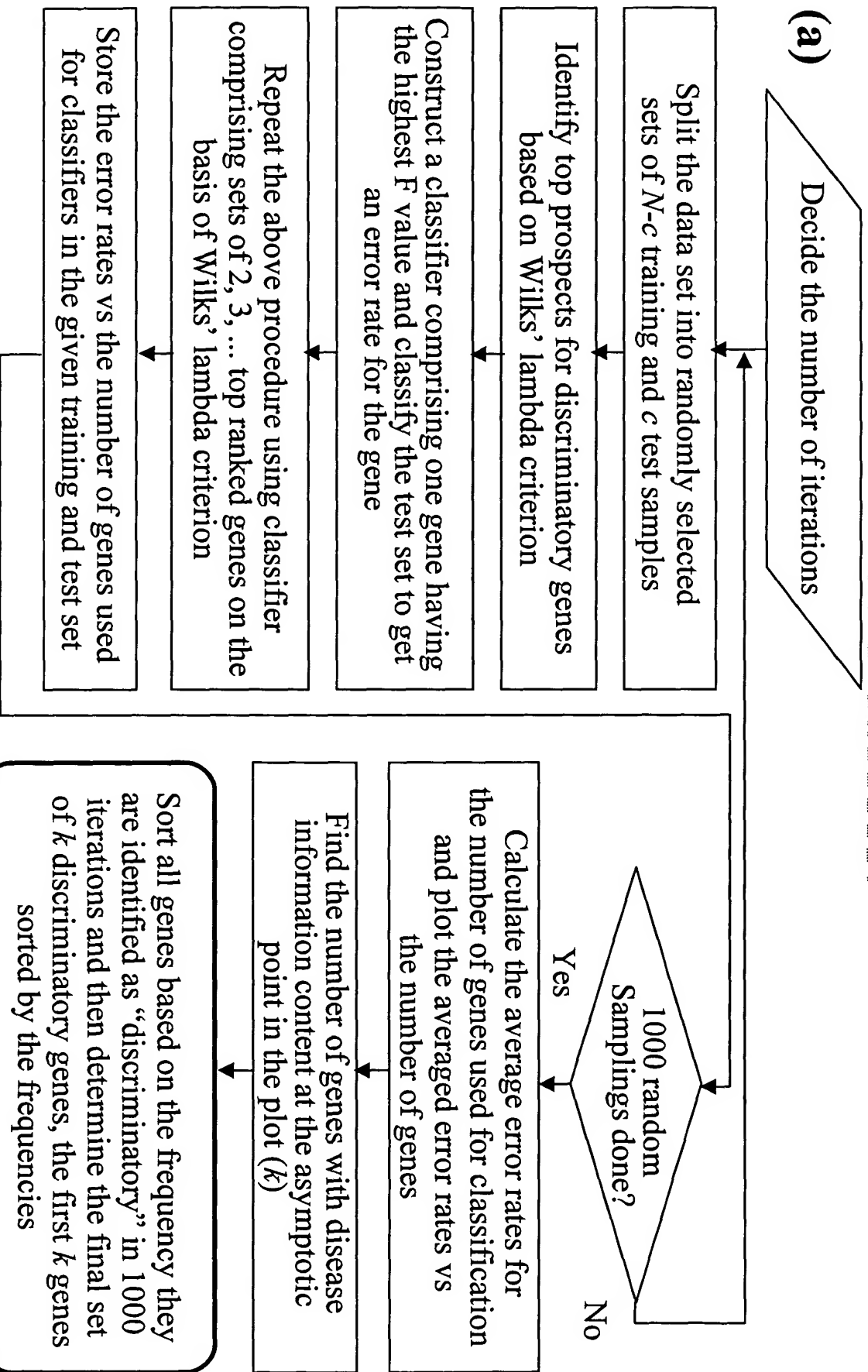


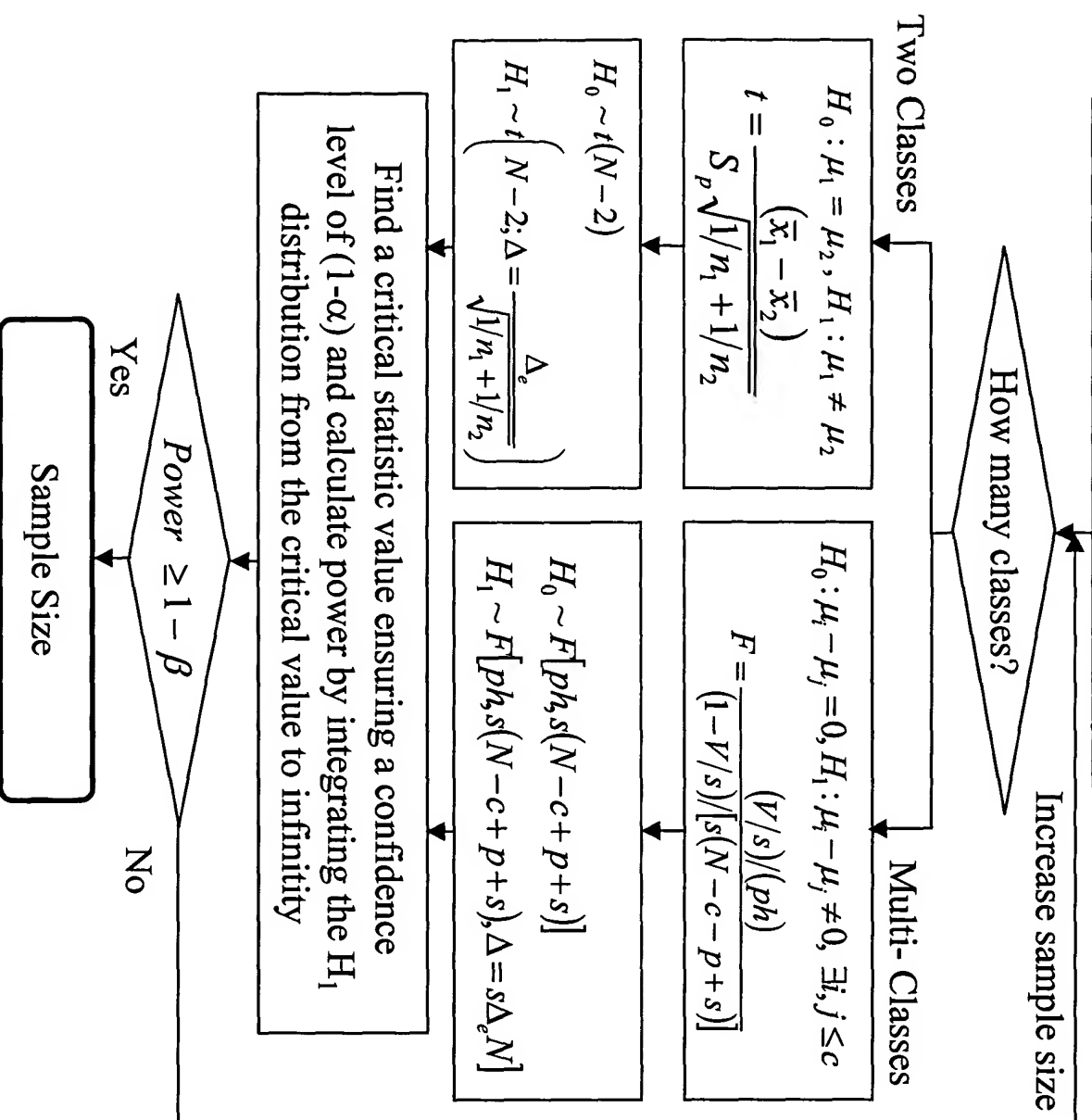
Figure 4: a) Projection of the expression phenotypes of cultures of *Synechocystis* sp. PCC 6803 to a CDA-defined space. This photosynthetic bacterium was grown under conditions shown in b) and the expression levels of 87 genes were measured by a DNA microarray at 29 time points spanning the entire course of the experiment. Of the 87 genes, 27 were identified as most discriminating of the four classes defined by the four different light conditions and their expression levels were projected to the CDA-defined space. It can be seen that the four phenotypic classes are clearly identified in the 3-dimensional CDA projection space. c) The first CV shows the largest discrimination power separating all the groups, discriminating clearly Phase III from the others. The second CV separates Phase IV from Phases I and II, while the third CV is necessary to separate Phase I from II.

(a)



**Figure 5a.** Leave one out cross-validation (LOOCV) algorithm, where  $N$  is the total number of samples and  $c$  is the number of classes, so that one sample from each class is included in the test.

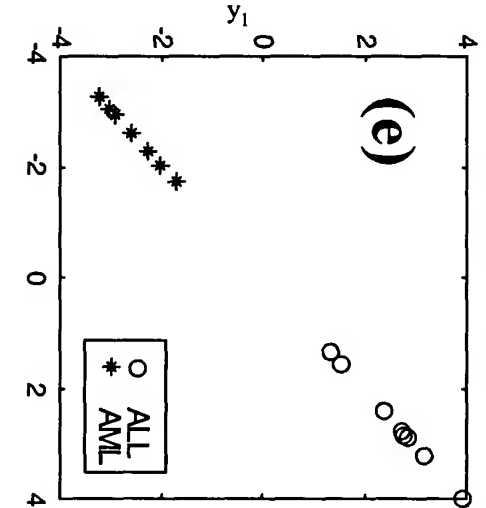
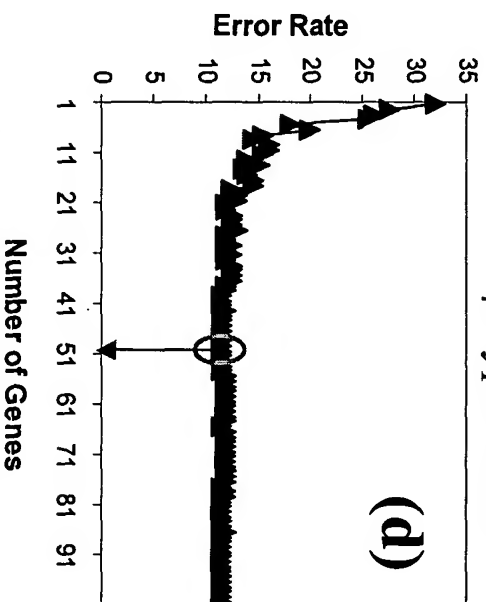
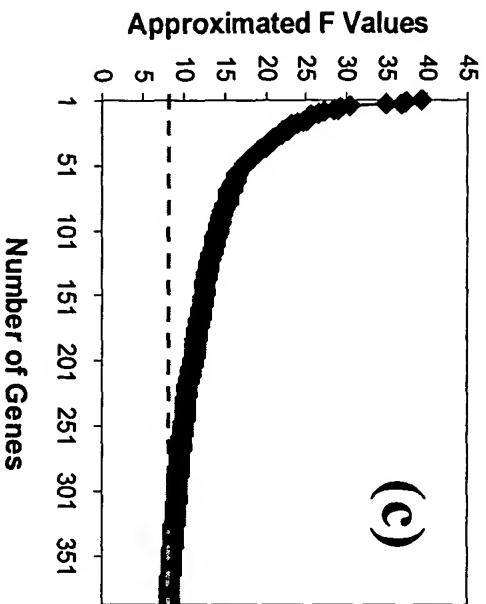
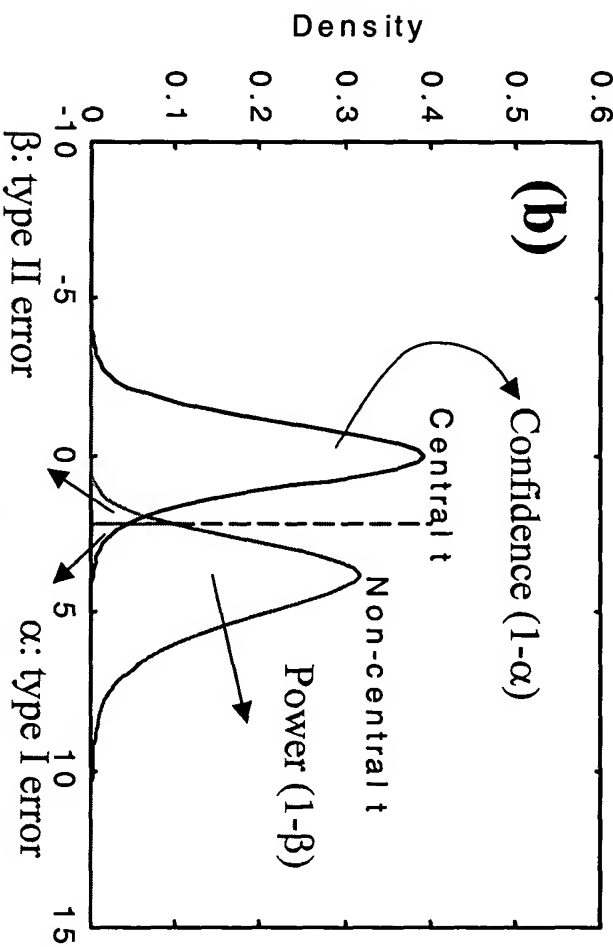
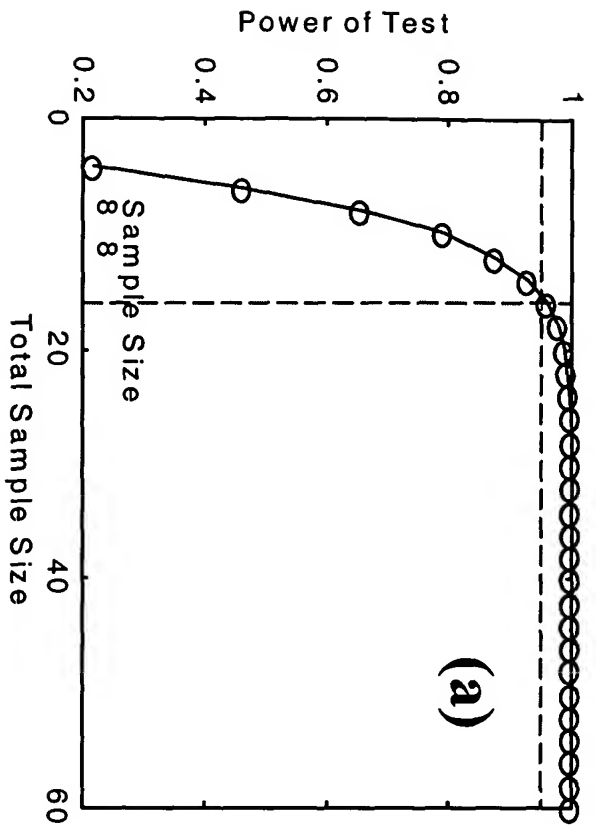
(b)



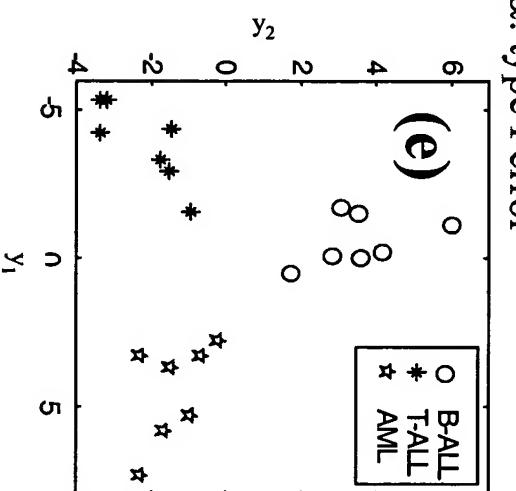
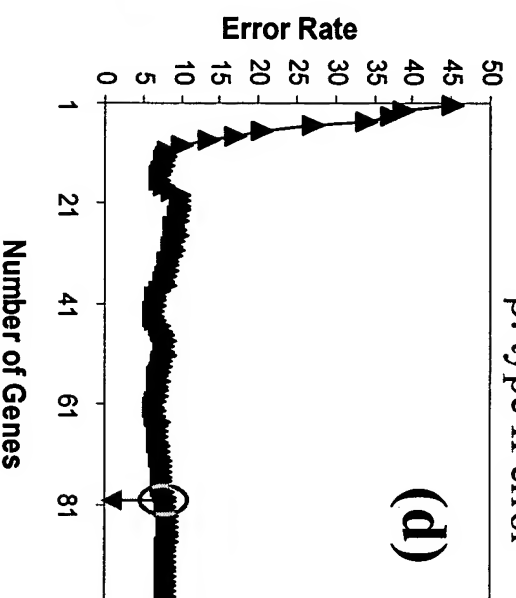
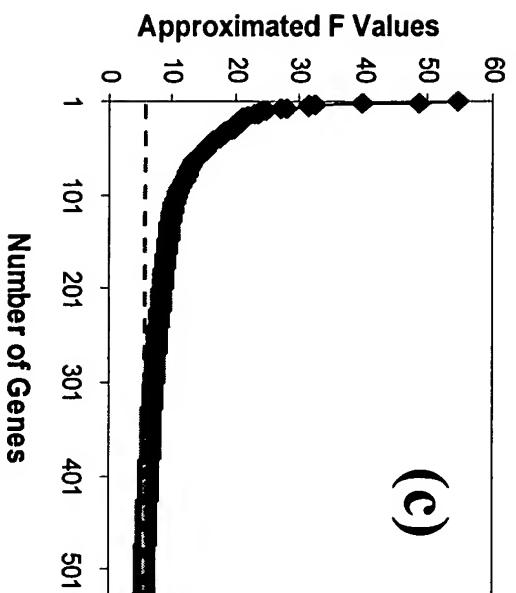
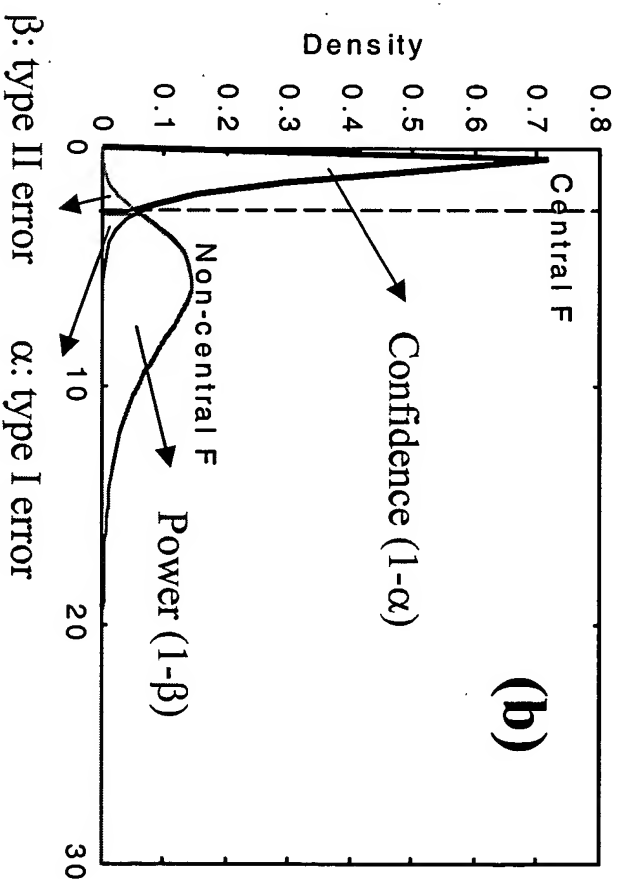
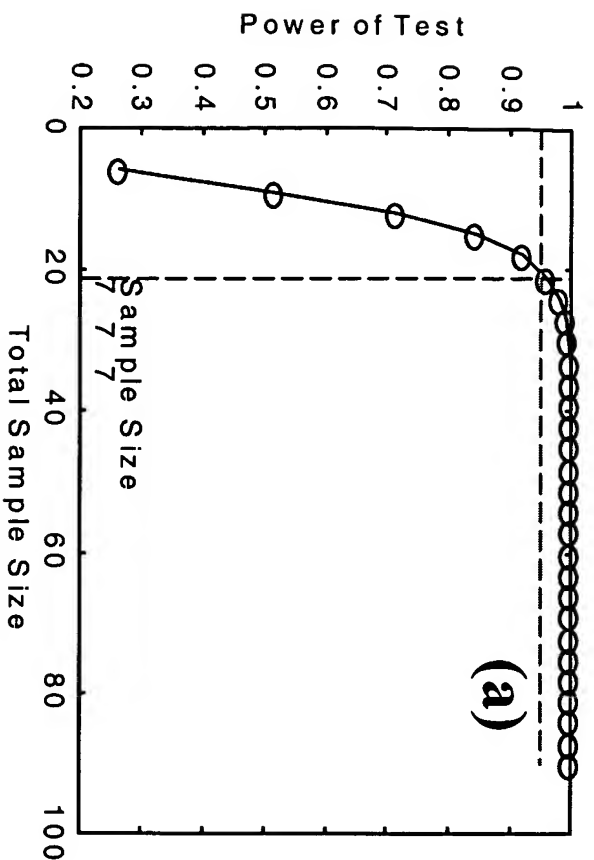
Design of hypothesis test

Define distributions  
of H<sub>0</sub> and H<sub>1</sub>

Figure 5b. Power analysis algorithm for determination of the minimum sample size.

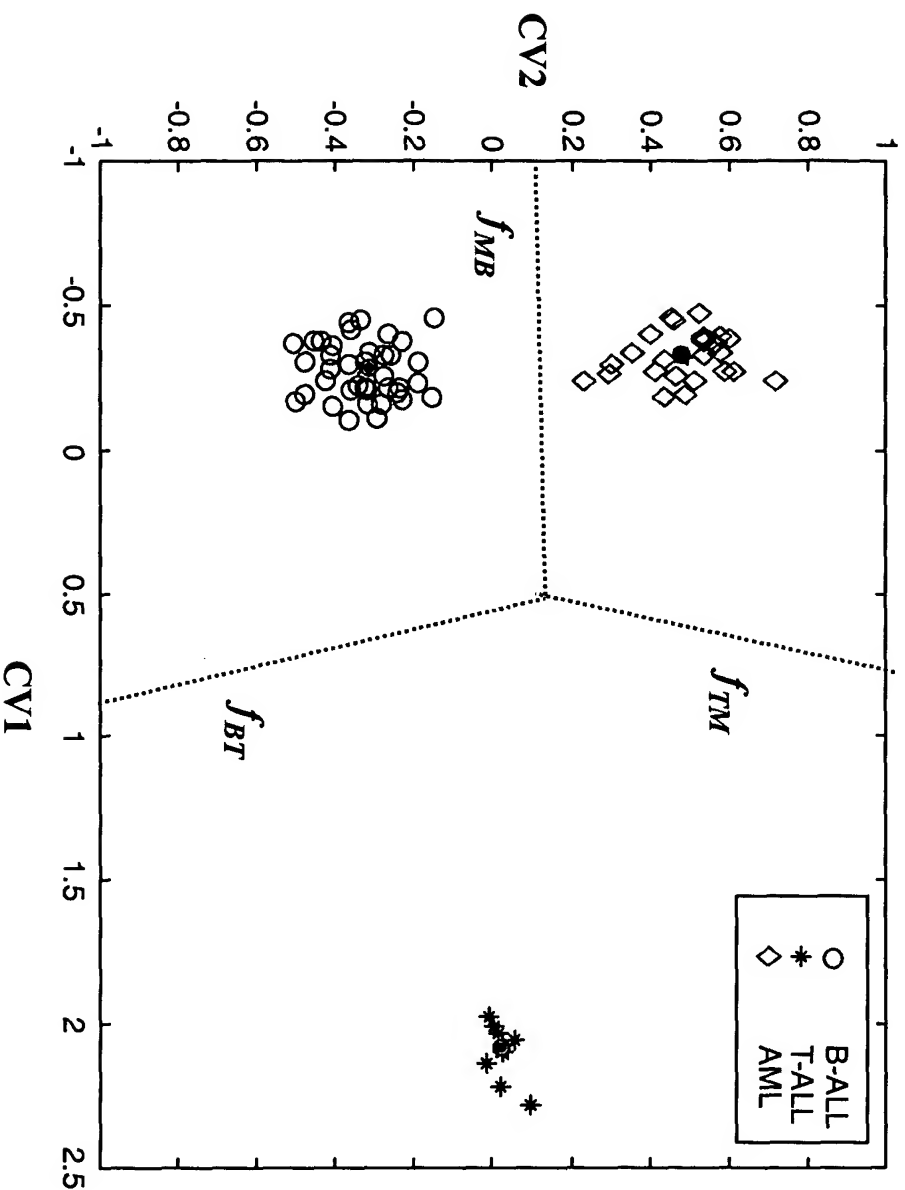


**Figure 6.** determination of minimum sample size for two-class (ALL, AML) distinction, selection of discriminatory genes with the estimated sample sizes of two classes, and FDA projection. (a) Power plot versus sample size showing how to determine the sample size required for two class distinction (8 from each class). (b) The distributions of  $H_0$  and  $H_1$  for the determined sample size. (c) Univariate F statistic values of the initial 388 discriminatory genes with a threshold ( $F_{0.01(1,18)} = 8.2854$ ) in randomly selected 8 ALL and 8 AML samples out of the entire data set. (d) Leave-one-out cross-validation applied to estimate the classification error rates and then to select the 50 most discriminatory genes with the same samples. (e) Separation of the 8 ALL and 8 AML samples in the two-dimensional FDA projection space defined discriminant axes of the 50 discriminatory genes.



**Figure 7.** determination of minimum sample size for the three-class (B-ALL, T-ALL, AML) distinction, selection of discriminatory genes with the estimated sample sizes of three classes, and FDA projection. (a) Power plot versus sample size showing how to determine the minimum sample size (7 from each class). (b) The distributions of  $H_0$  and  $H_1$  for the determined sample size. (c) Univariate F statistic values of the initial 527 discriminatory genes with a threshold ( $F_{0.01(2,26)} = 5.5263$ ) in randomly selected 7 B-ALL, 7 T-ALL and 7 AML samples out of the entire data set. (d) Leave-one-out cross-validation applied to estimate the classification error rates and then to select the 80 most discriminatory genes with the same samples. (e) Separation of the 7 B-ALL, 7 T-ALL and 7 AML samples in the two-dimensional FDA projection space defined by discriminant axes of the discriminatory 80 genes.





### Decision boundaries

$f_{BT}$ : between B-ALL and T-ALL

$f_{BM}$ : between T-ALL and AML

$f_{TB}$ : between AML and B-ALL

$$f_{mk} : y = V^T x$$

$$(\bar{y}_m - \bar{y}_k)^T y - \frac{1}{2}(\bar{y}_m - \bar{y}_k)^T (\bar{y}_m + \bar{y}_k) \geq 0$$

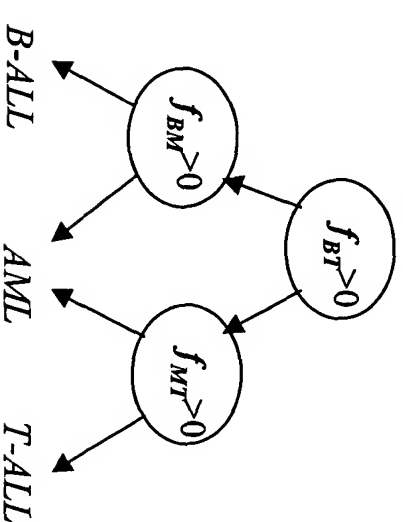


Figure 8: Physiological state domains are defined by a series of decision boundaries for the projected values in the reduced CV plane. Geometrical information can be inferred from this projection: for example, the covariance structure in B-ALL group is different from those in other groups. The decision boundary can be defined as the line perpendicular to the line joining two adjacent group means (each group mean is represented by a solid green circle). A classification tree can then be built for class prediction using those decision boundaries, as shown on the bottom right.

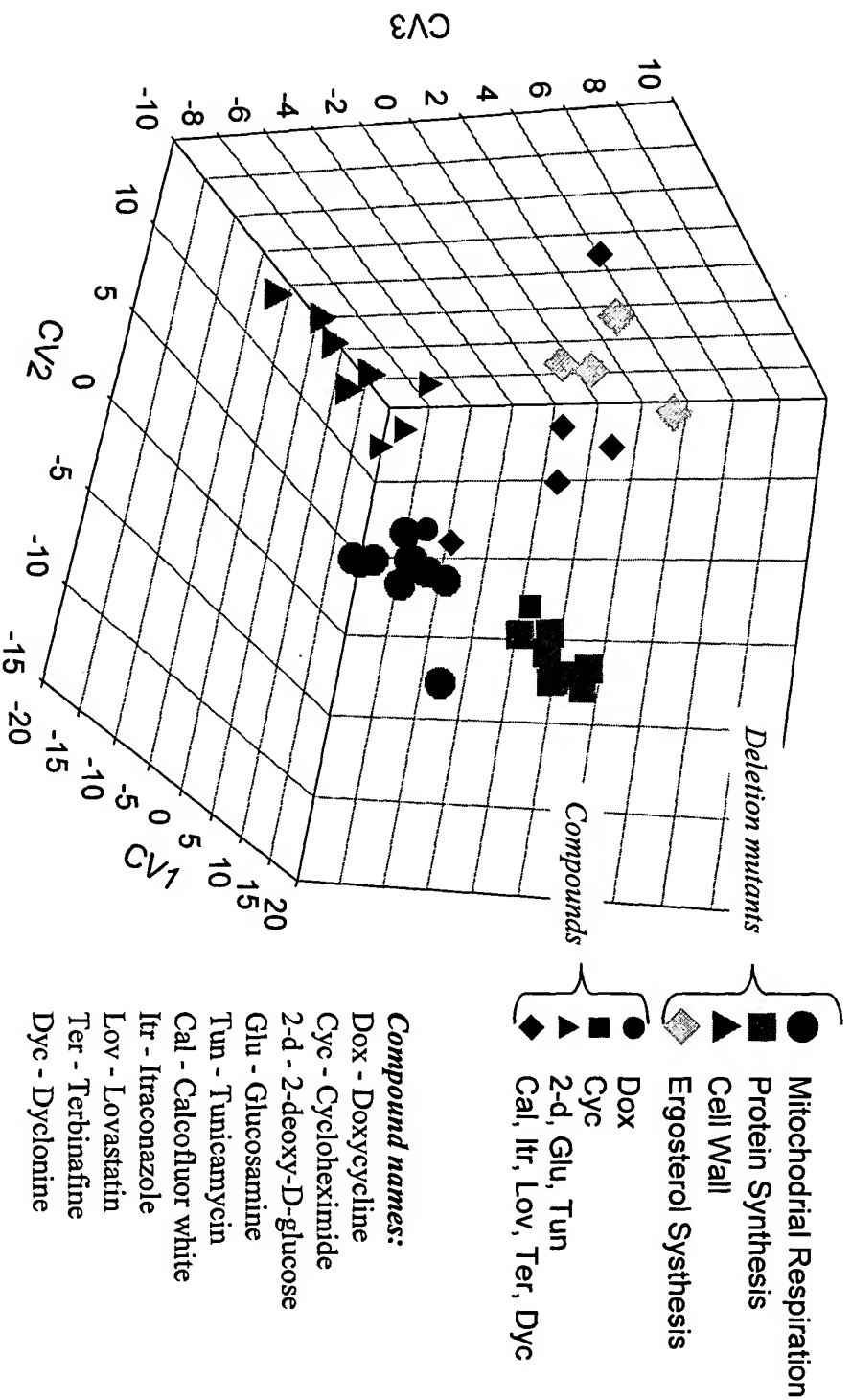
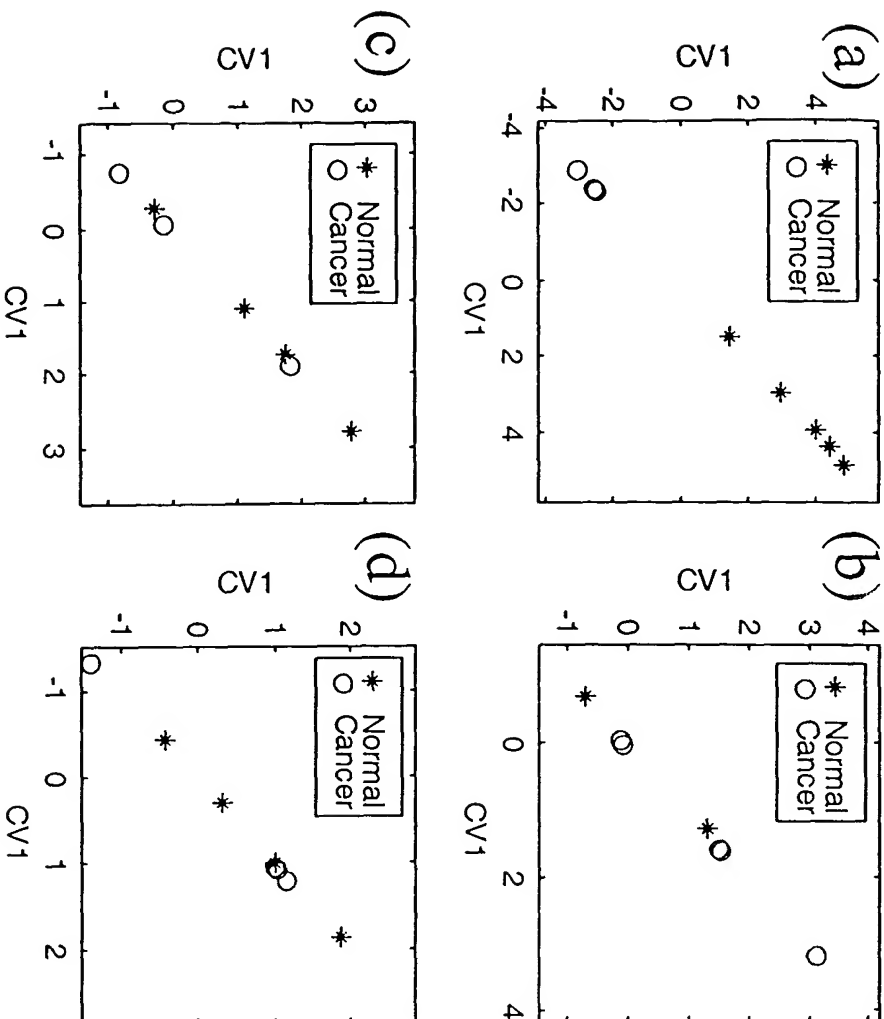


Figure 9: CDA projection of 27 yeast deletion mutant expression phenotype experiments grouped by the functionality of the eliminated gene. Four groups of related mutants have been distinguished using three CVs by projecting the expression levels of 200 of the most discriminating genes. The expression phenotypes obtained from the application of 10 chemical compounds to the wild-type yeast cultures are also projected into the CDA space defined by the mutants. The proximity in CDA space of these projections to those of the expression phenotype of the deletion mutant groups helps characterize the action of the compound on cell physiology. Note that one compound experiment (Cal) which appears incorrectly classified is actually in the center of the 3-D diagram, and not clearly associated with any of the groups shown. The classification suggested by the proximity of the projected phenotypes to the deletion mutants groups agrees with classification provided by Huges *et al.* (2000).



**Figure 1c.** CDA results using 45 genes. The projected value (CDA scores) into CV were calculated using linear combination of individual gene expressions,  $CV = v_1g_1 + v_2g_2 + \dots + v_{45}g_{45}$ . (a) Top 45 discriminatory genes, (b) 2000th-2045th genes, (c) 4000th-4045th genes, (d) 6000th-6045th genes. The overlap and variation in the groups increases when genes are chosen poorly.

Canonical Discriminant Analysis (CDA) filters out the most relevant discriminatory information from the data space by projecting the samples into a reduced space. This reduced space is defined by a set of Canonical Variables (CV), which are a linear combination of the discriminatory genes.

The following optimization problem is solved to compute each CV, denoted by  $v$ .

$$Max \frac{v^T W^{-1} B v}{v^T v}$$

$$W = \sum_{k=1}^c (X_k - \bar{X}_k 1^T) (X_k - \bar{X}_k 1^T)^T$$

$$T = \sum_{k=1}^c (X_k - \bar{X}_k 1^T) (X_k - \bar{X}_k 1^T)^T \quad B = T - W$$

Here  $W$  represents the within-group variance,  $B$  represents the between-group variance, and  $T$  represents the total variance. The solution is the matrix of eigenvectors ( $V$ ) of  $W^{-1}B$  whose corresponding eigenvalues ( $\Lambda$ ) represent the relative discriminatory power of each CV.

$$W^{-1}B = V \Lambda V^T$$

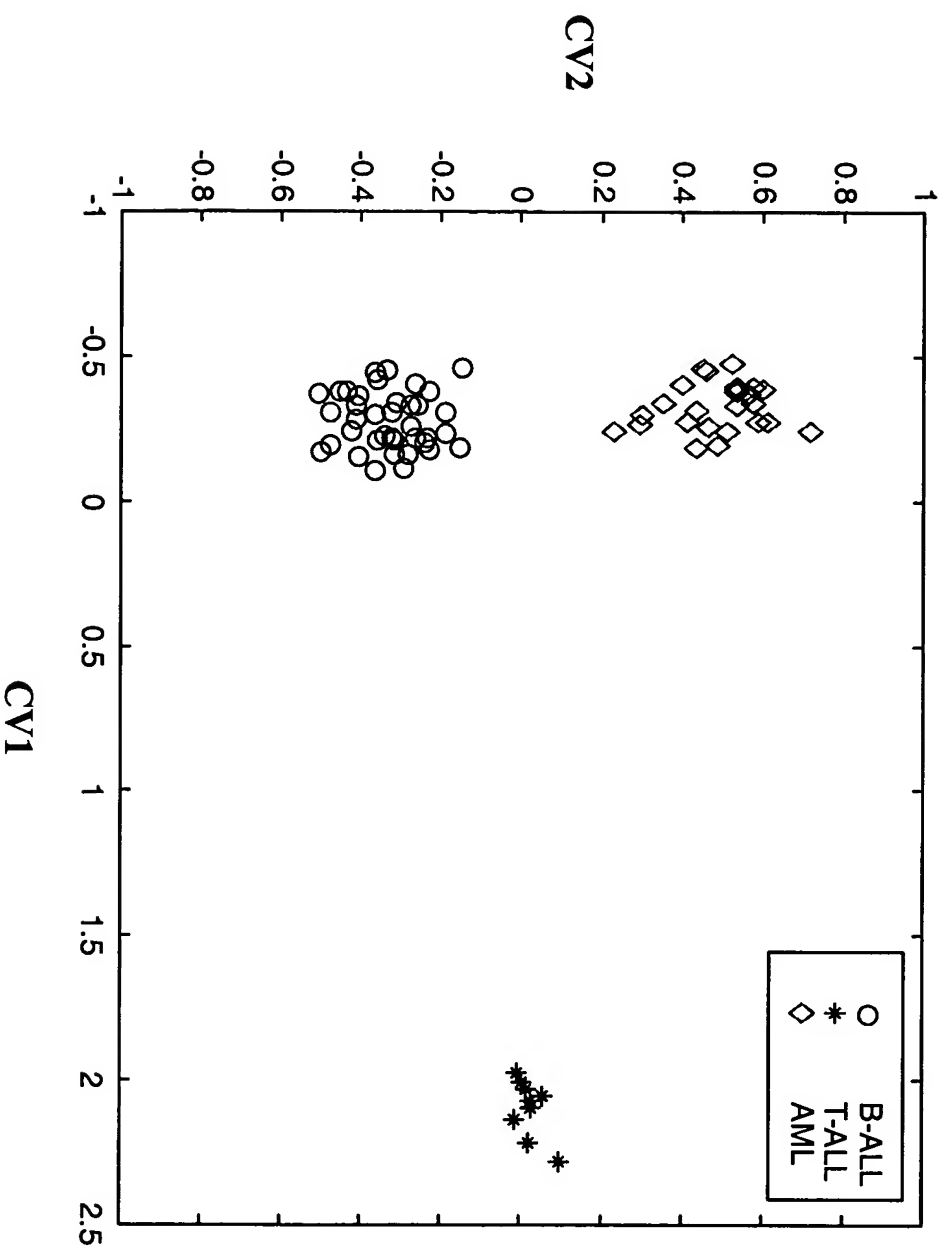


Figure 11: CDA projection of expression data obtained from patients with B-ALL, T-ALL, and AML. Projection of 30 discriminatory genes allow CDA to clearly separate the three classes of leukemia expression phenotype in a two dimensional CV plane. The first CV distinguishes the B-ALL group from T-ALL and AML. The second CV separates T-ALL group from AML to complete the group separation.

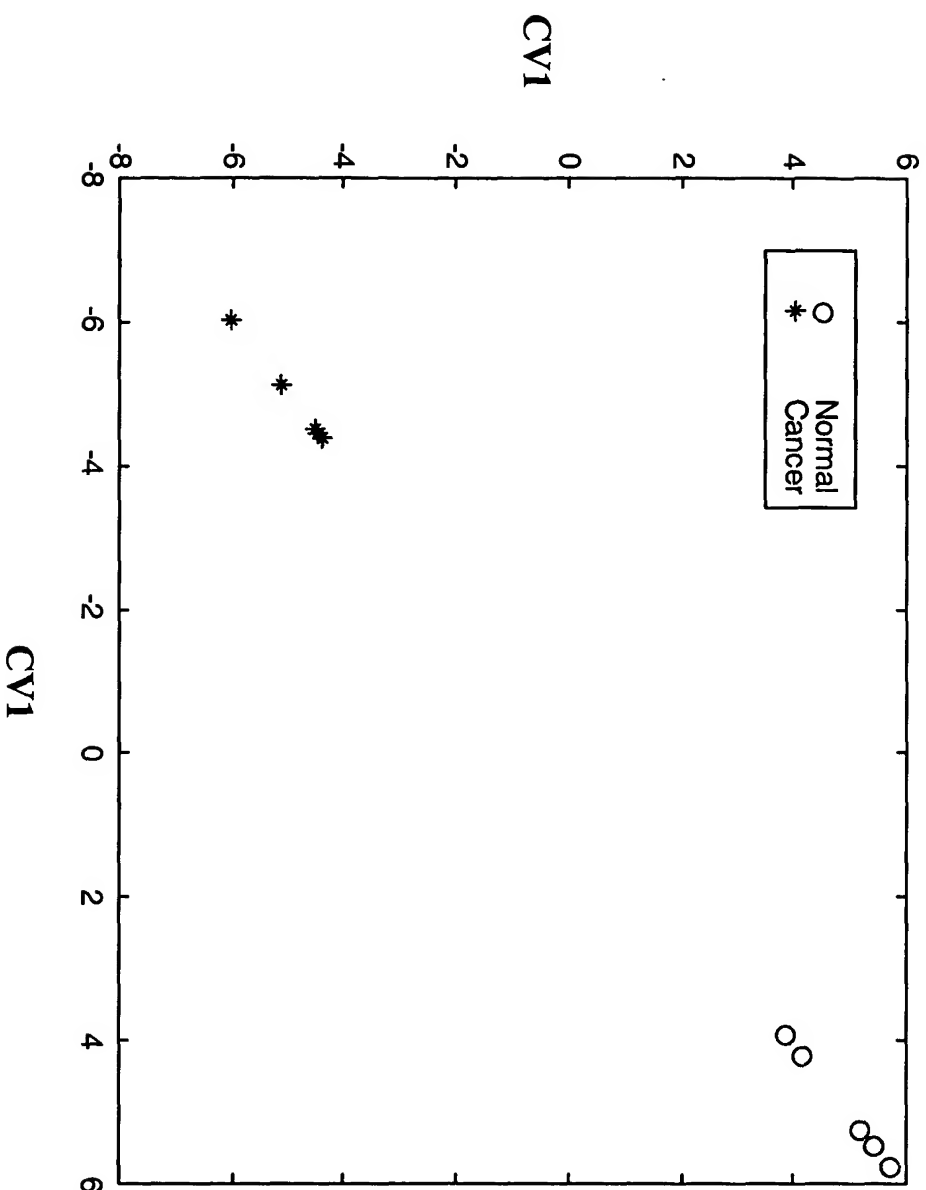
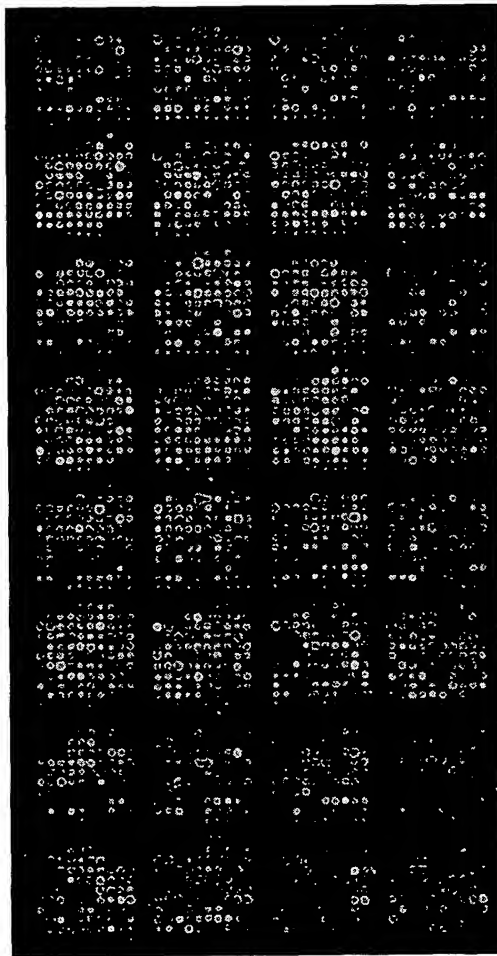


Figure 12: CDA projection of the expression phenotypes comprising 7070 genes measured in samples obtained from healthy individuals (5 samples) and patients with oral epithelium cancer (5 samples). 35 discriminatory genes out of 7070 total genes allow CDA to clearly separate the two groups in one dimensional CV line. The CV value distinguishes normal tissues from cancerous tissues.

-69)-



THE UNIVERSITY OF CHICAGO

Figure 13

bioRxiv preprint doi: <https://doi.org/10.1101/170170>; this version posted July 1, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

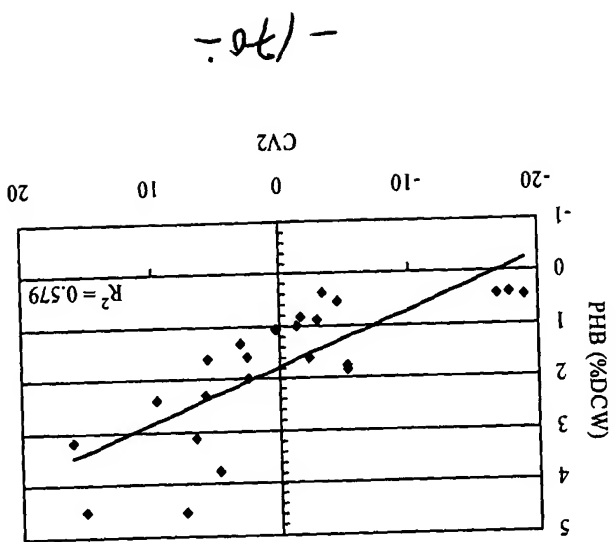
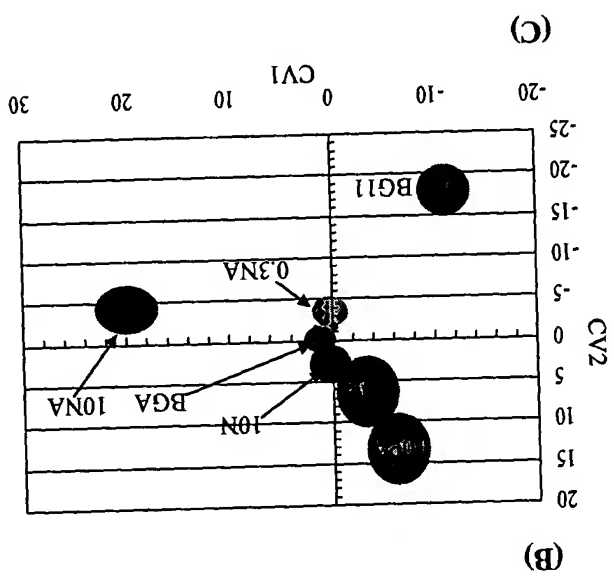
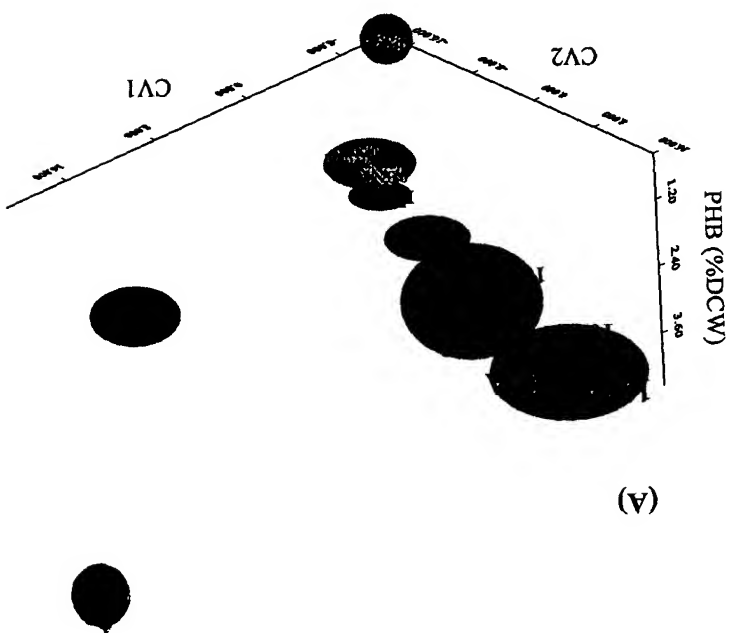


Figure 14

bioRxiv preprint doi: <https://doi.org/10.1101/111111>; this version posted January 1, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

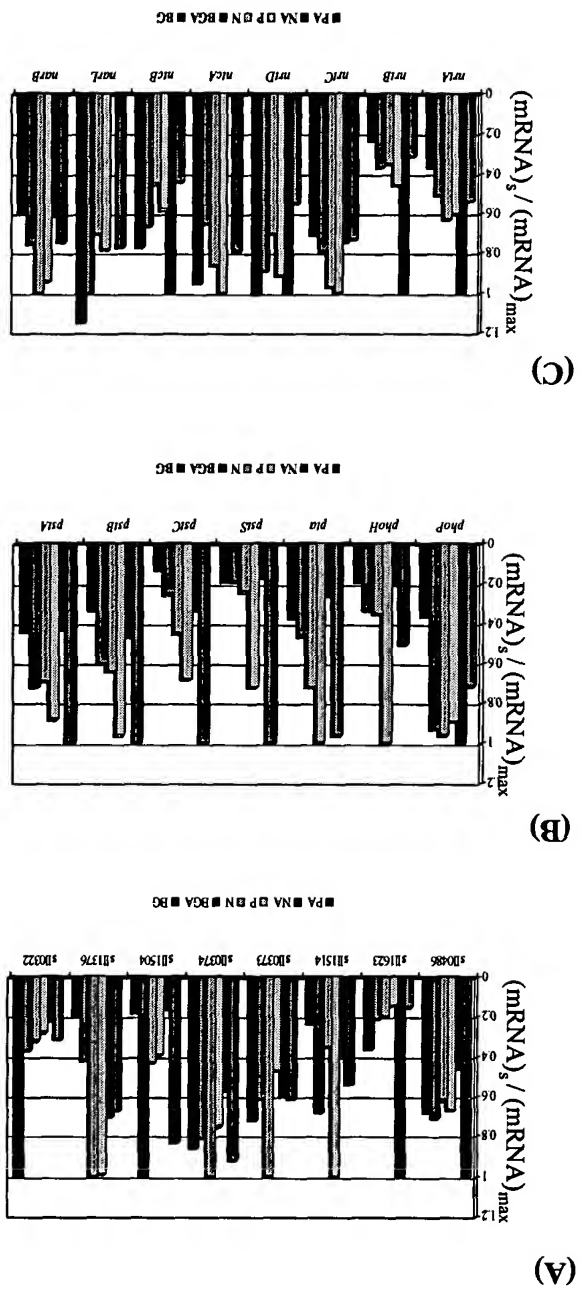


Figure 15

- (17) -



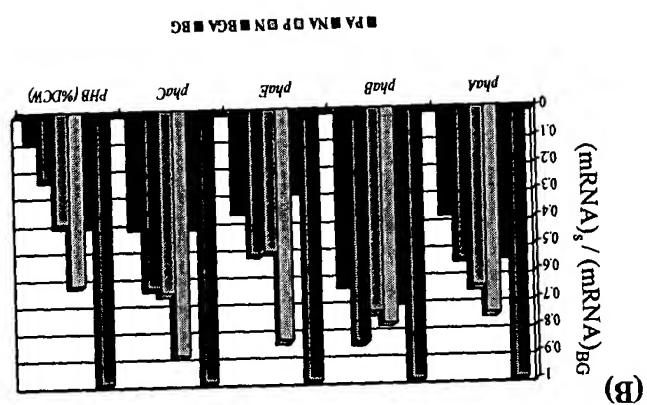
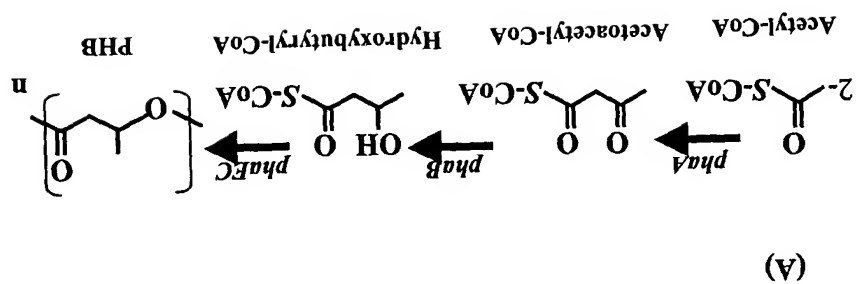


Figure 16